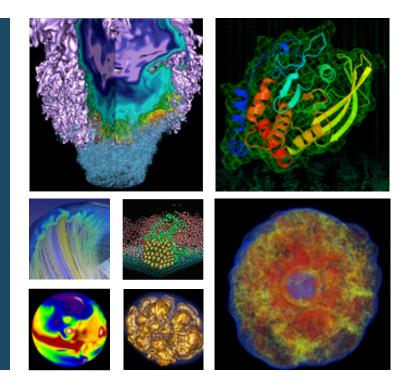# A Year in the Life of a Parallel File System



**Glenn K. Lockwood**, Shane Snyder, Teng Wang, Suren Byna, Philip Carns, Nicholas J. Wright

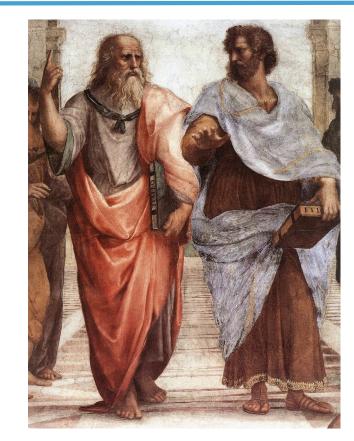**November 15, 2018**

# Why was my job's I/O slow?



Socrates (left) and Plato (right) contemplating I/O performance in *The School of Athens* by Raphael. 1511.

# Why was my job's I/O slow?

1. You are doing something wrong
2. Another job/system task is competing with you
3. The storage system is degraded

# Why was my job's I/O slow?

1. You are doing something wrong
2. Another job/system task is competing with you
3. The storage system is degraded

**Most frustrating**

**Least studied**

# Our holistic approach to I/O variation

1. Measure performance variation over a year on large-scale production HPC systems
2. Collect telemetry from across the entire system
3. <u>Quantitatively</u> describe why I/O varies so much

# 1. Observing variation in the wild

- **Probe I/O performance daily**
  - Jobs scaled to achieve >80% peak fs performance
  - 45 – 300 sec per probe
- **Run in diverse production environments**
  - Two DOE HPC facilities (ALCF, NERSC)
  - Three large-scale systems (Mira, Edison, Cori)
  - Two parallel file system implementations (GPFS, Lustre)
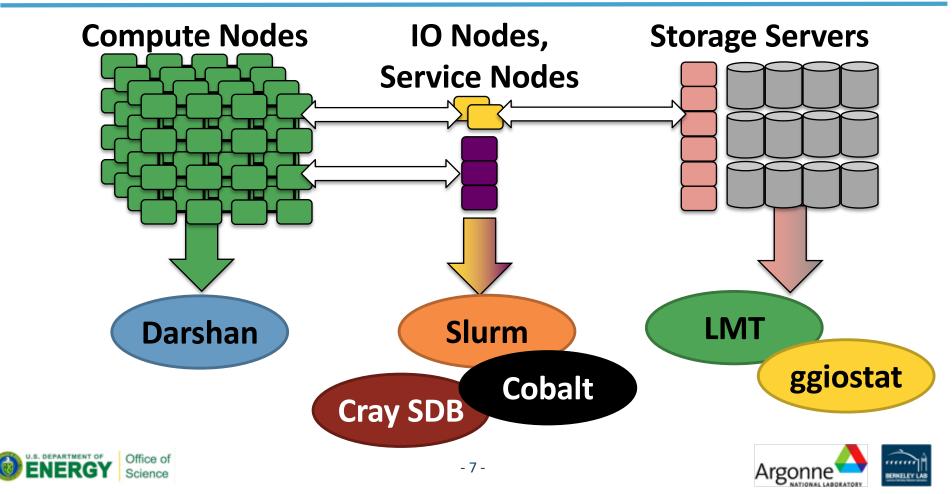  - Five file systems (Mira gpfs1, Edison lustre[1-3], Cori lustre1)

| App I/O Transfer Size | Shared File | File Per Process |
|---|---|---|
| O(1 MiB) | IOR | IOR |
| O(100 MiB) | VPIC BD-CATS | HACC |

# 2. Collecting diverse data for holistic analysis

# Year-long I/O performance dataset



- **366 days of testing**

- **11,986 jobs run**

- **220 metrics measured per job**
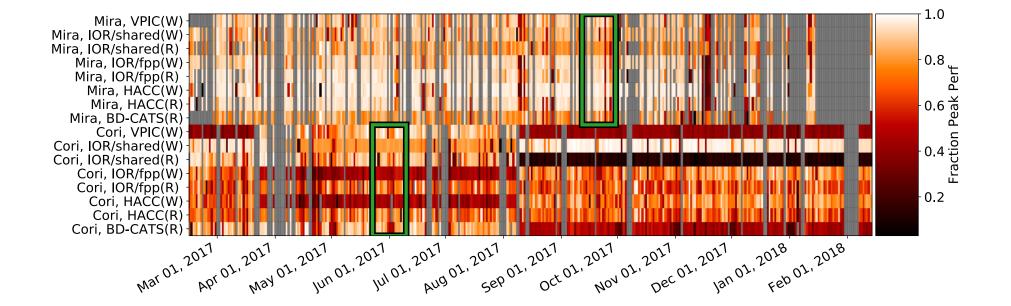  - some derived or degenerate
  - sometimes undefined

**…and not very insightful at a glance**
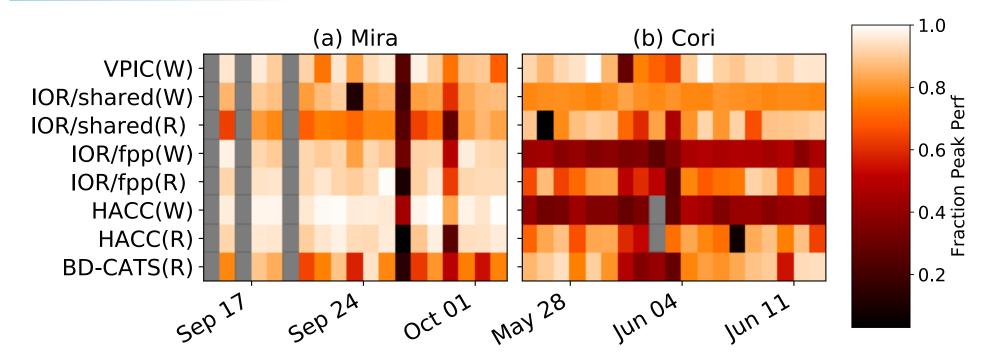
# I/O performance variation in production
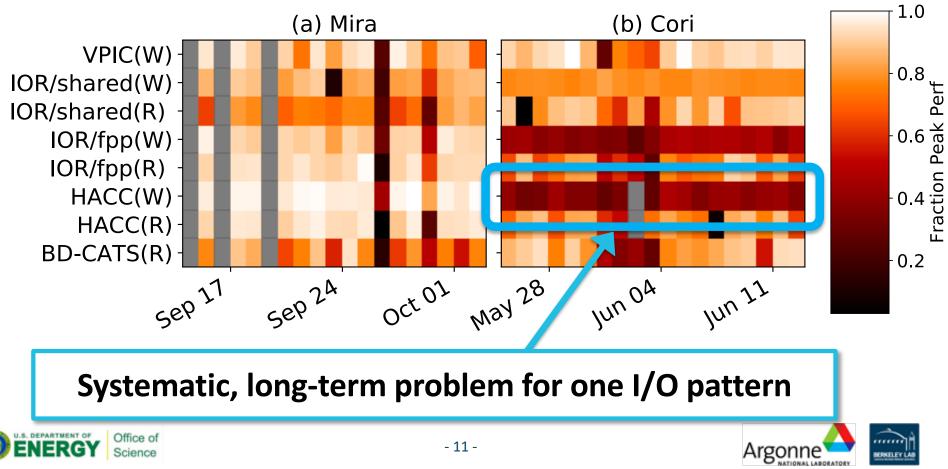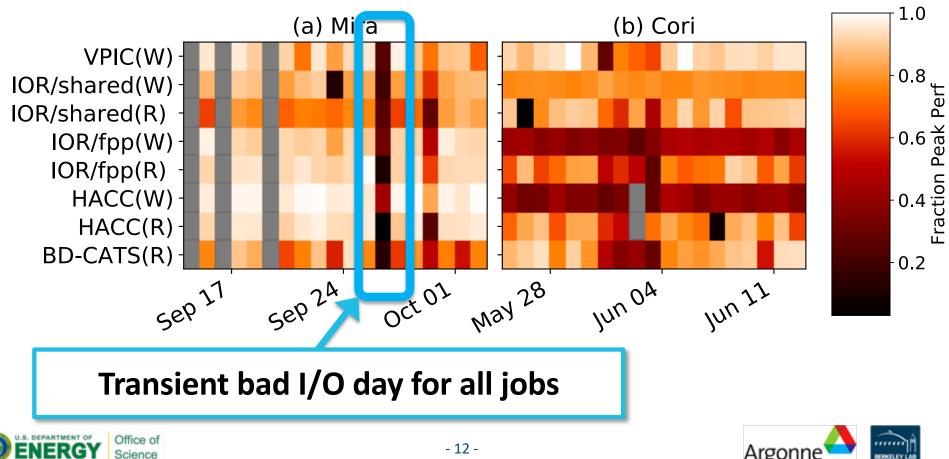
# Two flavors of I/O performance variation



(a) Mira      (b) Cori

# Performance varies over the long term



(a) Mira      (b) Cori

**Systematic, long-term problem for one I/O pattern**

# Performance varies over the short term



**Transient bad I/O day for all jobs**

# Performance also experiences transient losses



(a) Mira     (b) Cori

**Transient I/O problems**

# Again: Why was my job's I/O so slow?

- **Could be:**
  - Long-term systematic problems
  - Short-term transient problems
- **The next questions:**
  - What causes long-term, systematic problems?
  - What causes short-term transient problems?
- **Our approach:**
  - Separate problems over these two time scales
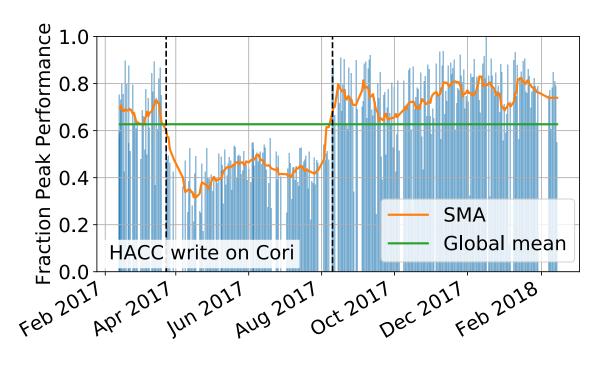  - Independently classify causes of longer-term and shorter-term variation

# Separating short-term from long-term

- **Goal: Numerically distinguish time-dependent variation**

- <u>S</u>imple <u>m</u>oving <u>a</u>verages (SMAs) from financial market technical analysis

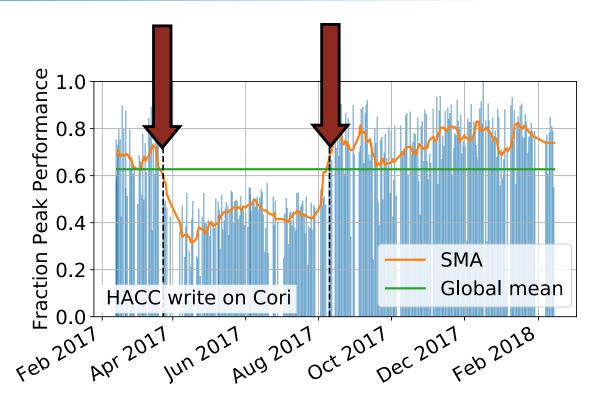- Where short-term average performance diverges from overall average

# Quantitatively bound long-term problems

- **Goal: Numerically distinguish time-dependent variation**

- <u>S</u>imple <u>m</u>oving <u>a</u>verages (SMAs) from financial market technical analysis

- Where short-term average performance diverges from overall average

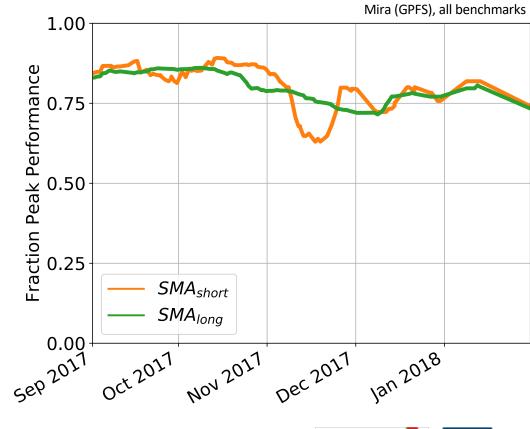- Example: Bug in a specific file system client version

**Goal: Contextualize transient variation happening during long-term variation**

- Two SMAs at different time windows
  (e.g., 14 days and 49 days)
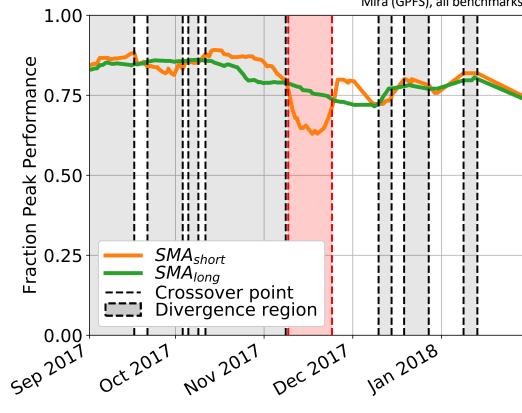
Mira (GPFS), all benchmarks

# Separating short-term from long-term variation

**Goal: Contextualize transient variation happening during long-term variation**

- Two SMAs at different time windows
  (e.g., 14 days and 49 days)

- <u>Crossover points</u> indicate short behavior == long behavior

- <u>Divergence regions</u> where short behavior diverges from long behavior
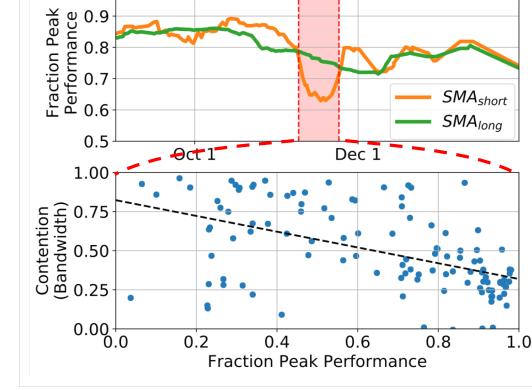
Mira (GPFS), all benchmarks



Legend:
- $SMA_{short}$ (orange)
- $SMA_{long}$ (green)
- Crossover point (dashed)
- Divergence region

Y-axis: Fraction Peak Performance (0.00, 0.25, 0.50, 0.75, 1.00)
X-axis: Sep 2017, Oct 2017, Nov 2017, Dec 2017, Jan 2018

# What causes divergence regions?

- Capitalize on widely ranging performance (and all 219 other metrics)
- Correlate performance in this region with other metrics
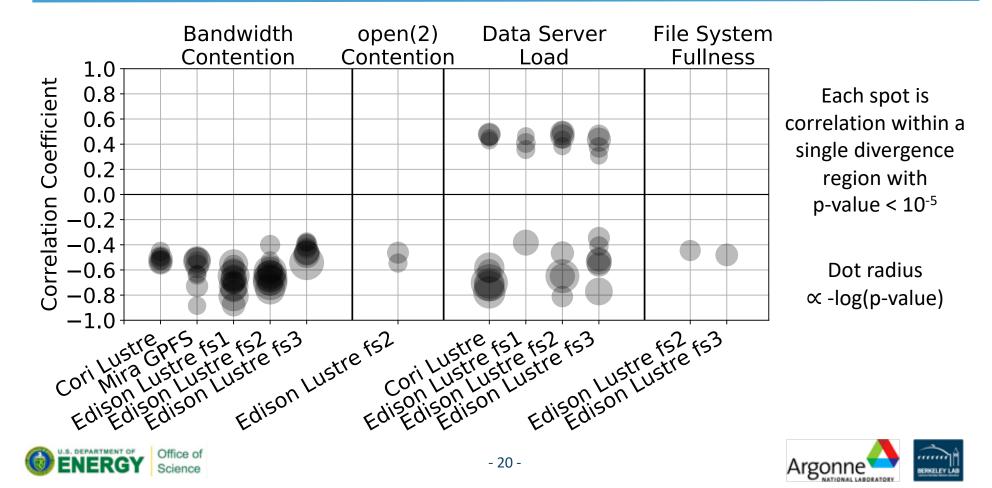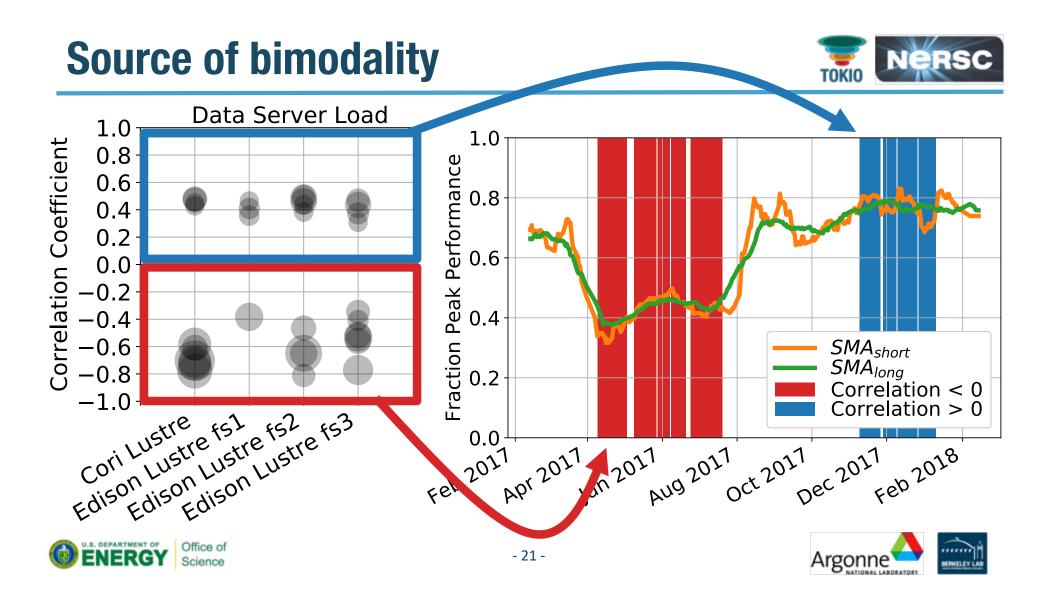  - Bandwidth contention
  - IOPS contention
  - Data server CPU load
  - ...



Mira (GPFS), all benchmarks

# What causes short-term variation over a year?

# Source of bimodality

# Identifying sources of transient variation

Mira (GPFS), all benchmarks

- Partitioning allows us to classify short-term performance variation
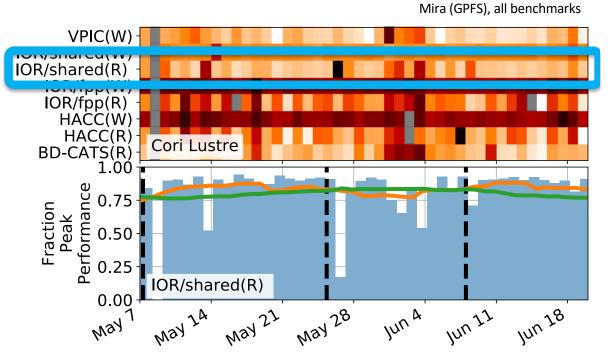
- Can't correlate truly transient variation though

# Identifying sources of transient variation

- Confidently classifying transients is statistically impossible

- Classifying in aggregate *is* possible!

- If we observe a possible relationship…
    - One time?  Maybe coincidence
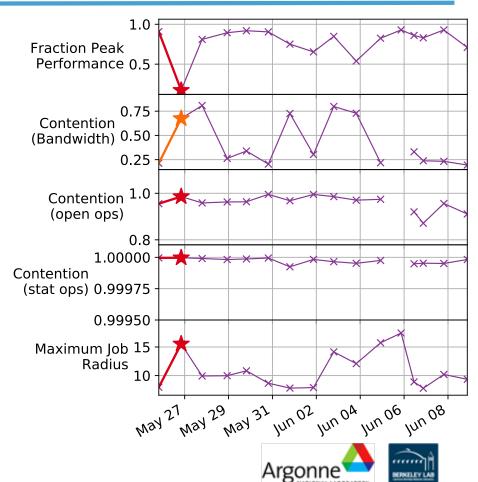    - Many times? Maybe *not* a coincidence



Mira (GPFS), all benchmarks

1. Identify jobs affected by transient issues

2. Define divergence regions

3. Classify jobs based on region, calculate p-values

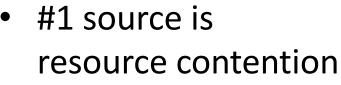4. Repeat for all transients and, calculate aggregate p-values

# Sources of transient variation in practice
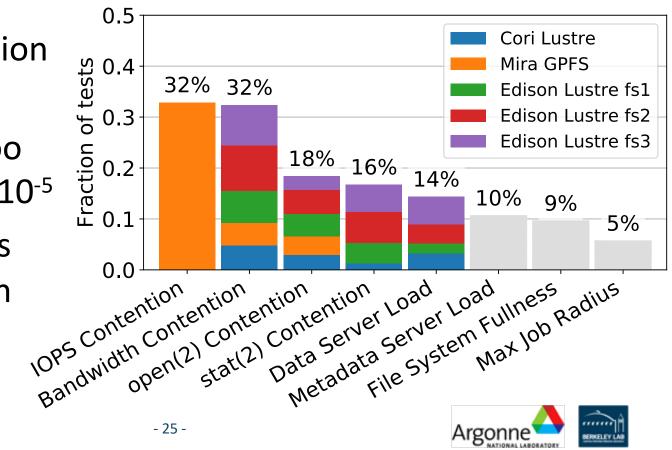
- #1 source is resource contention

- Other factors implicated but too rare to meet $p < 10^{-5}$

- 16% of anomalies defy classification

# Overall findings

- **Baseline performance and variability change over time**
  - Patches & updates
  - Sustained bandwidth contention from scientific campaigns

- **Partitioning performance in time yields more insight**
  - Can classify short-term and transient variation
  - Quantifies effects of contention and suggests avenues for system architecture optimization

- **We can learn things from other fields of study**

# Try this at home!

**Reproducibility (code + year-long dataset):**

https://www.nersc.gov/research-and-development/tokio/a-year-in-the-life-of-a-parallel-file-system/ (or see the paper appendix)

**pytokio Framework:**

https://github.com/nersc/pytokio